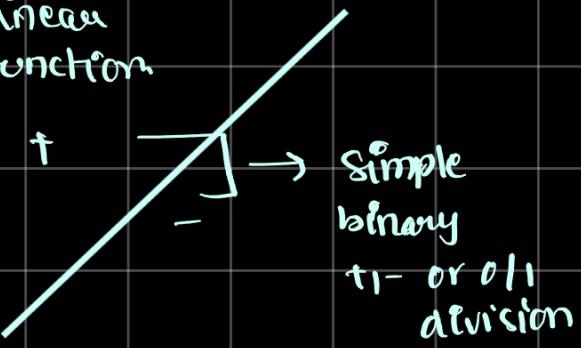
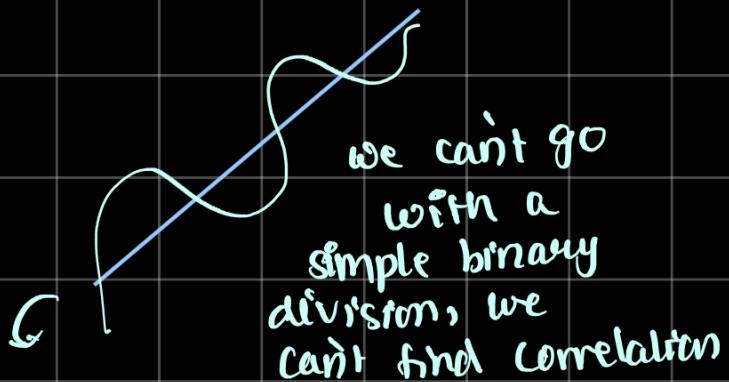


Linear function.



Non-linear curve



thus we find individual relations

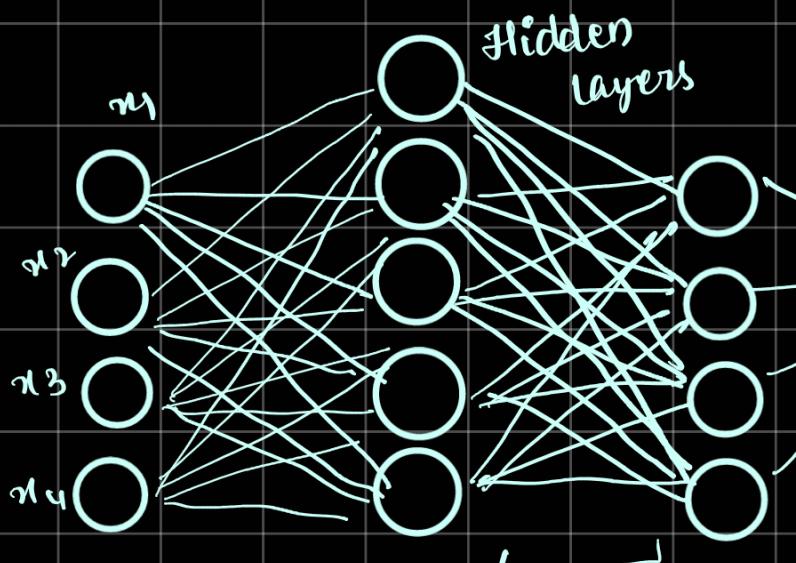
$$\left. \begin{array}{l} \text{ie } x_1 \rightarrow y_1 \\ x_2 \rightarrow y_2 \\ \vdots \\ x_n \rightarrow y_n \end{array} \right\} \rightarrow \text{finding a correlation}$$

then fitting it with a function

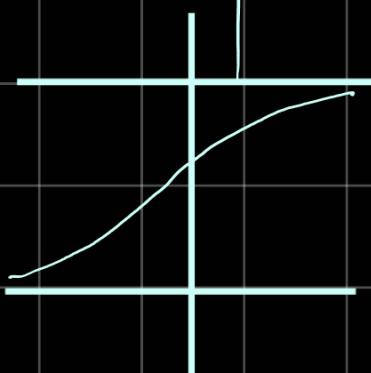
feature engineering

domain knowledge \rightarrow for a machine to predict this features

we use deep learning



adds non-linearity



logistic regression

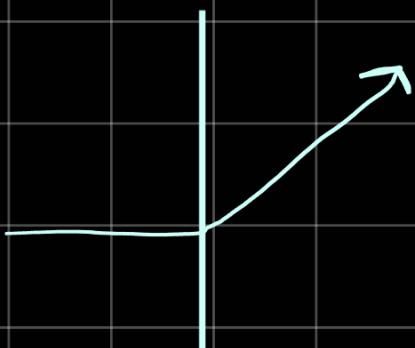
$$y = \sigma(\theta_0x_0 + \theta_1x_1 + \dots + \theta_nx_n)$$

$$G \sigma(\theta_0x_0 + \theta_1x_1 + \dots + \theta_nx_n)$$

σ sigmoid graph

Alternative activation function to sigmoid function is

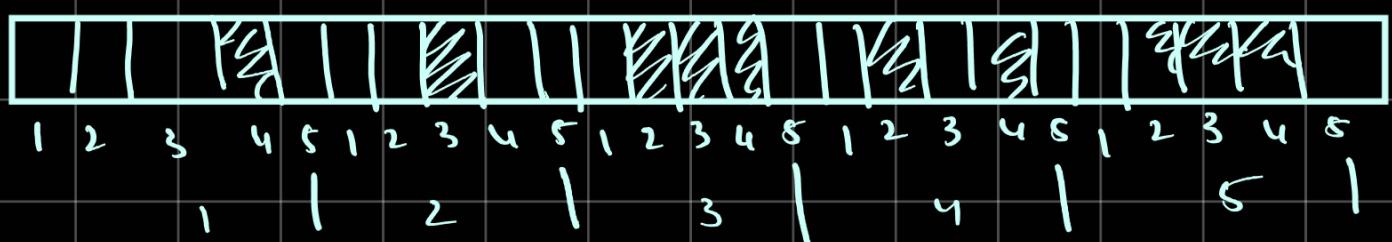
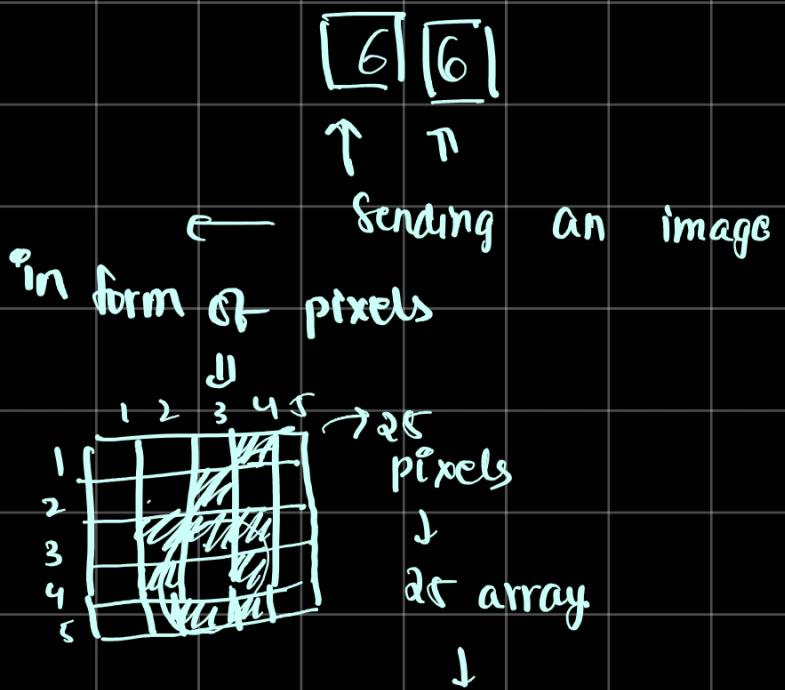
ReLU
= ↗ (Popular)

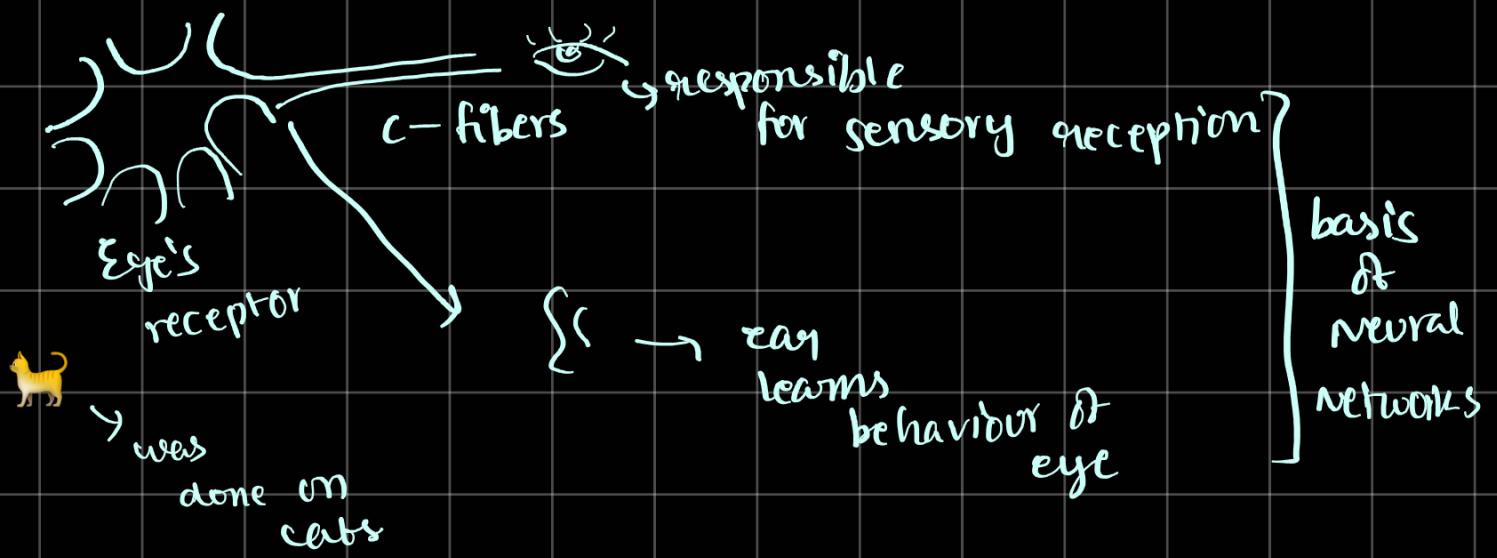


→ partially linear function

→ discontinuous but
Lasso applicable.

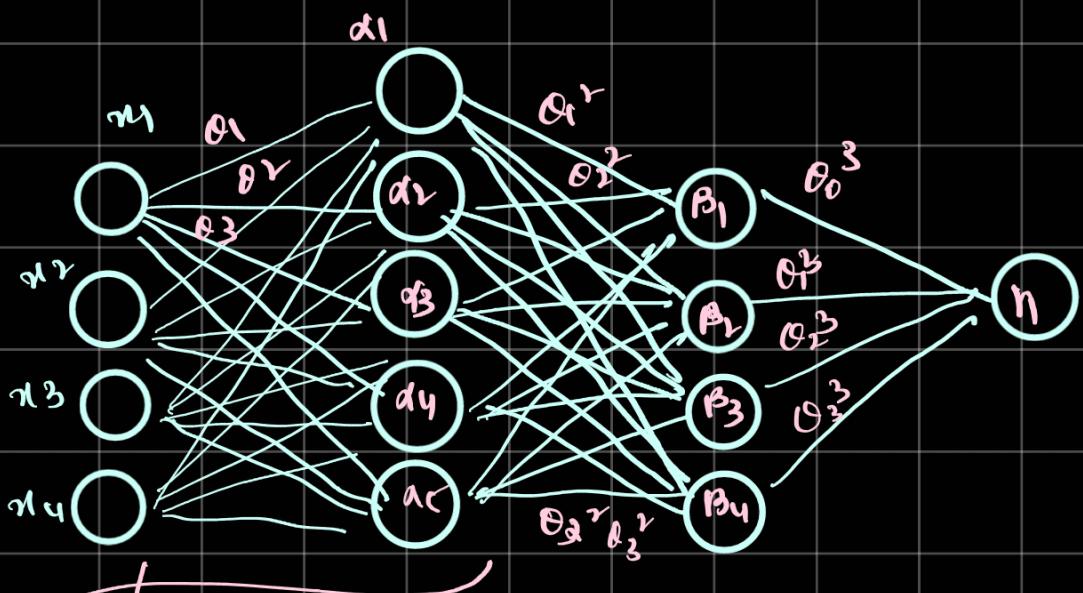
↳ captures features from local to global to determine the universal features of an element in a neural network





→ if $m \rightarrow$ inputs
 $n \rightarrow$ hidden networks
 $\hookrightarrow m \times n$

|
 |
 $n \rightarrow$ inputs
 $p \rightarrow$ hidden networks
 $\hookrightarrow n \times p$



$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4 \dots \theta_n x_n$$

1st Layer i input

$$\theta_0 + \theta_1 d_1 + \theta_2 d_2 + \theta_3 d_3 + \dots + \theta_k d_k$$

Hidden layer 1

$$\theta_0 + \theta_1 a_1 + \theta_2 a_2 + \theta_3 a_3 + \theta_4 a_4 + \dots + \theta_u a_u$$

Hidden layer k

Loss function J

$$J = -\gamma \ln(\gamma') - (1-\gamma) \ln(1-\gamma')$$

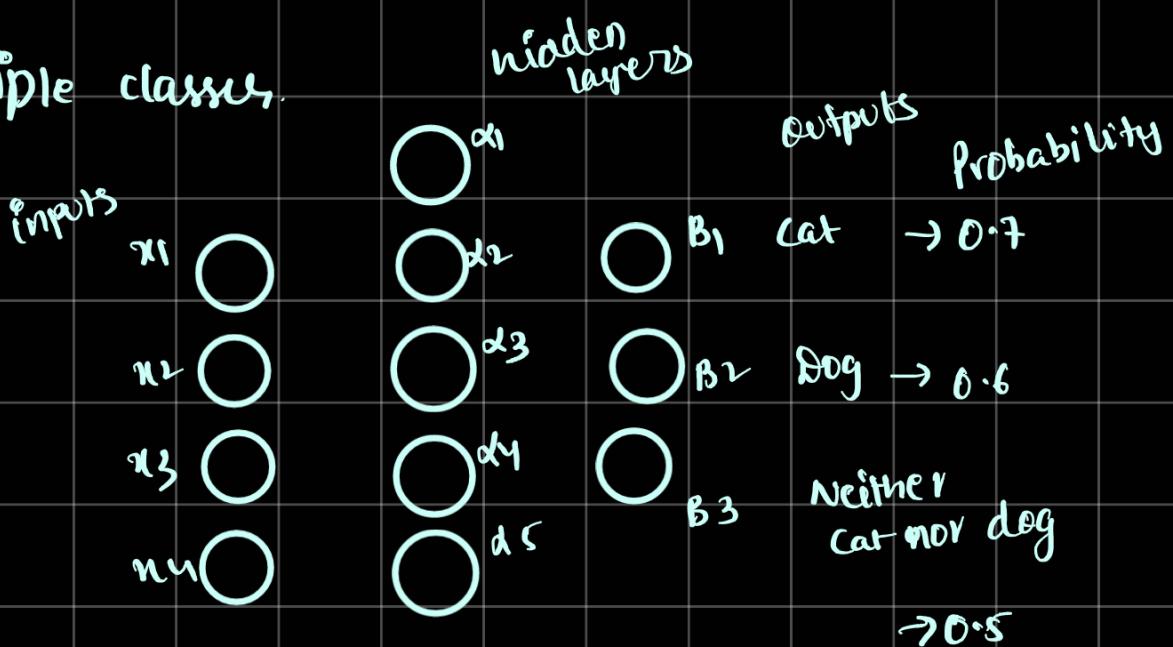
$$\frac{\partial J}{\partial \theta_i} = \frac{\partial J}{\partial B_1} \times \frac{\partial B_1}{\partial \theta_i} \rightarrow \text{using chain rule}$$

↓
each hidden unit focuses on a local feature ; one class at a time

→ forward passes in neural networks act as "tensors"
not as loops

$$; \begin{bmatrix} x_1 & x_2 & x_3 & x_4 \\ x_1^2 & x_2^2 & x_3^2 & x_4^2 \\ x_1^3 & x_2^3 & x_3^3 & x_4^3 \end{bmatrix}_{3 \times 4} \times \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \end{bmatrix}_{4 \times 1} \rightarrow \text{Simple tensor multiplication.}$$

→ Multiple classes.



Soft max

↳ when probabilities are really close, it scales the differences so that the differences become significant

Say probability of cat = 0.7

$$\text{Soft max} = \frac{e^{0.7}}{e^{0.7} + e^{0.6} + e^{0.5}}$$

↓ ↓
Dog neither cat | Dog

(linear layer)

ReLU ↳ Rectified Linear Unit → non activation function

↳ The output of the 1st hidden layer is sent to a ReLU function to apply the function



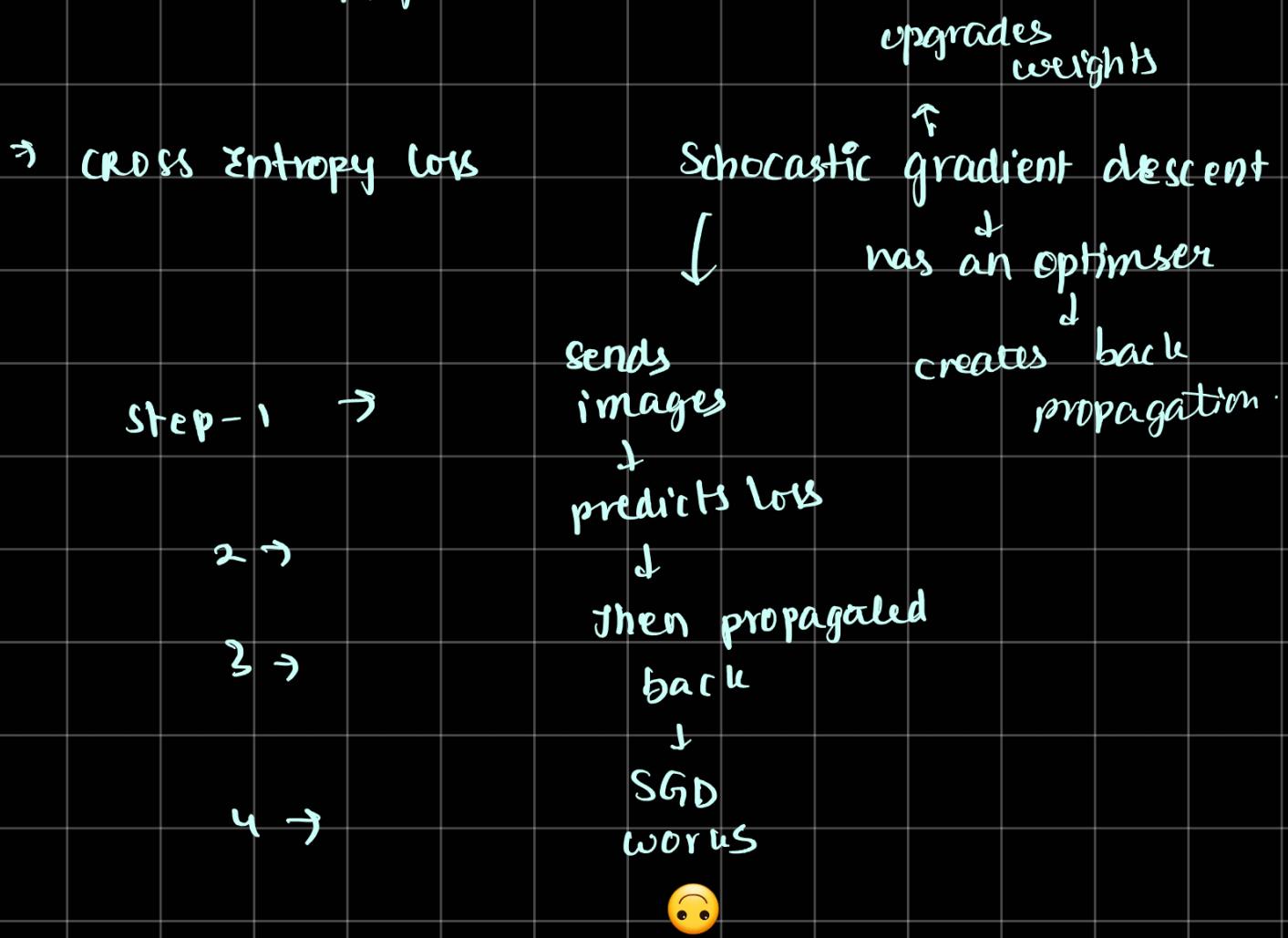
$$\text{ReLU}(x) = \begin{cases} 0 & x \leq 0 \\ x & x > 0 \end{cases}$$

Label is the one which has the highest probability

→ Depth of network is n-1 cause input is not counted as a layer.

Forward computation & Backward loss propagation
↓
rates get adjusted
↳ The Q's are the rates.

Hyperparameters - The specific parameters that we give while training our model.



Accuracy vs and loss vs for each epoch

Too many epoch → over fitting

constant value → model converged.

